# Vision and Language
## LXMLS 2024
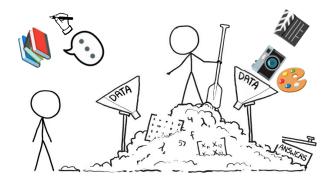
### Desmond Elliott

Department of Computer Science

University of Copenhagen

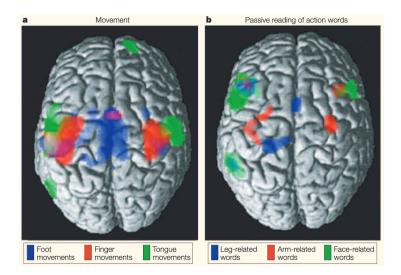Slides: https://elliottd.github.io/vlprimer/

# Working Definition

Multimodal models jointly processes information from two or more input modalities, e.g. images and text, speech and video, etc.
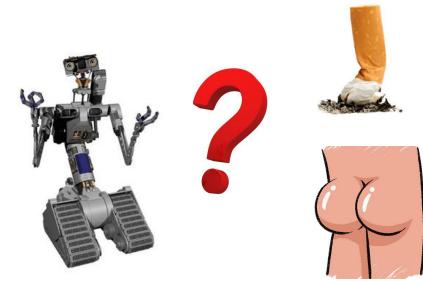
# Why Multimodality?

- Humans ground conceptual knowledge in modality processing systems in the brain

- Evidence that grounding activates similar brain regions for different input modalities

Barsalou et al. (2003). Grounding conceptual knowledge in modality-specific systems. Trends in cognitive sciences, 7(2):84–91.
Pulvermüller. (2005). Brain mechanisms linking language and action. Nature reviews neuroscience, 6(7), 576-582.
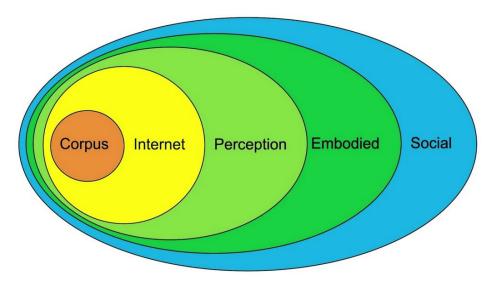
3

# Multimodality reduces ambiguity

# You Cannot Learn Language From

- The radio without grounding
  *(lack perception)*

- The television without actions
  *(lack embodiment)*

- *W*ithout interacting with others
  *(lack social)*



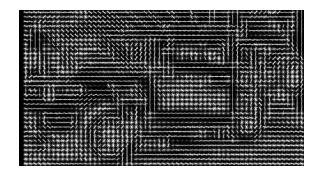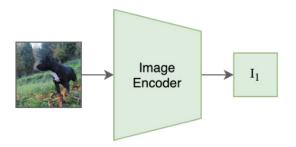Bisk et al. EMNLP 2020.. Experience Grounds Language.

# (At Least) Five Major Areas

- **Representation**: how to convert raw inputs into a usable format

- **Translation**: transform from one modality to another

- **Alignment**: predict relationships between elements across modalities

- **Fusion**: join features from modalities to support prediction

- **Co-learning**: transferring knowledge from one modality to another

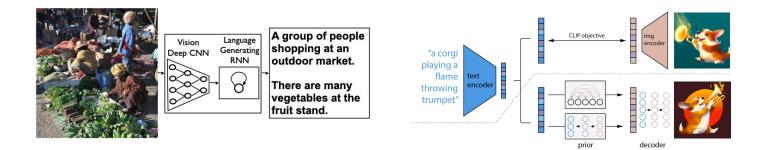Baltrušaitis, Ahuja, and Morency, (2018). Multimodal machine learning: A survey and taxonomy. IEEE PAMI, 41(2), 423-443.

# Representation

- Great deal of work over the last decade, from HOG features in the early 2000s to CLIP features in the 2020s.

Dalal & Triggs. CVPR 2005. Histograms of oriented gradients for human detection.
Radford et al. ICML 2021. Learning transferable visual models from natural language supervision.

# Translation

- Explosion of end-to-end neural network models since the mid 2010s

Vinyals et al. (2015). Show and tell: A neural image caption generator. CVPR.
Ramesh et al. (2022). Hierarchical Text-Conditional Image Generation with CLIP Latents. arXiv.

# Alignment

- Important for self-supervised learning and also for phrase grounding



Positive narration candidates $\mathcal{P}$
MIL-NCE positive contribution
Sampled negative narrations $\mathcal{N}$
Standard MIL positive contribution

A pink elephant: 1.00

Miech et al. (2020). End-to-end learning of visual representations from uncurated instructional videos. CVPR.
Kamath et al. (2021). MDETR-modulated detection for end-to-end multi-modal understanding. ICCV.
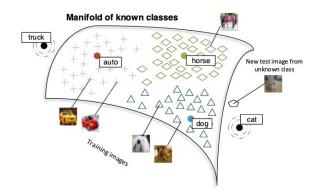
# Fusion

- Early work studied the differences between early and late fusion.
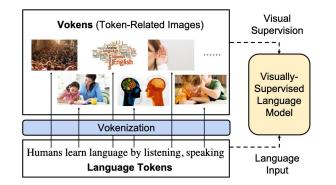- Multi-head self-attention now provides model-based fusion.



Chen and Jin (2016). Multi-modal conditional attention fusion for dimensional emotion prediction. MM.
Lu et al. (2019). ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *NeurIPS.*

# Co-learning

- Zero-shot transfer across modalities, or using visual grounding to improve language models on text-only tasks.

Socher et al. (2013). Zero-shot learning through cross-modal transfer. NeurIPS.
Tan & Bansal. (2020). Vokenization: Improving Language Understanding with Contextualized, Visual-Grounded Supervision. EMNLP

# Roadmap

Part 1

1. Datasets and Tasks for Multimodal Learning

   📚 Visually Grounded Reasoning across Languages and Cultures

2. Data Representation
3. Modelling Techniques

   📚 Retrieval-Augmentation for Image Captioning, Sequential Multimodal Compositional Generalization

Part 2

4. Understanding Multimodal Models

5. Future Directions

   📚 Language Modelling with Pixels

# 1. Datasets and Tasks for Multimodal Learning

# Two Types of Dataset

- **General-purpose**: visual data with descriptive annotations
    - Conceptual Captions
    - LAION-2/5B
    - Speech-COCO

Blue Beach Umbrellas, Point Of Rocks, Crescent Beach, Siesta Key - Spiral Notebook

- **Task-specific**: visual data with e.g. classification labels
    - Image / Video Captioning
    - Visual Question Answering
    - Visually Grounded Reasoning

What color is the cat's leash?
purple          red

# Many Types of Tasks

- Sequence generation
  - Image captioning, video captioning
    visual storytelling, image generation

$$P(x|v)$$

- Classification
  - VQA, Visually-grounded Reasoning

$$P(y|x,v)$$

- Ranking and Alignment
  - Image↔Text Retrieval
    Referring Expression Localization

$$\mathbf{Distance}(x,v)$$

# COCO

- Used both a **general-purpose** and **task-specific** dataset

- Images covering 80 *common* objects in *context* with multiple human-authored captions.

- Object segmentation data too!



some sheep walking in the middle of a road
a herd of sheep with green markings walking down the road
a herd of sheep walking down a street next to a lush green grass covered hillside.
sheared sheep on roadway taken from vehicle, with green hillside in background.
a flock of freshly sheered sheep in the road.

Chen et al. (2015). Microsoft COCO captions: Data collection and evaluation server. *arXiv.*

16

# Conceptual Captions

- Used for pretraining

- 3/12M images released with *normalized* English captions.

- Normalization is not public.

- Due to *linkrot*, much less data is currently available.

Sharma et al. (2018). Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. ACL.
Changpinyo et al. (2021). Conceptual 12M: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. CVPR.
Download your own: https://github.com/igorbrigadir/DownloadConceptualCaptions
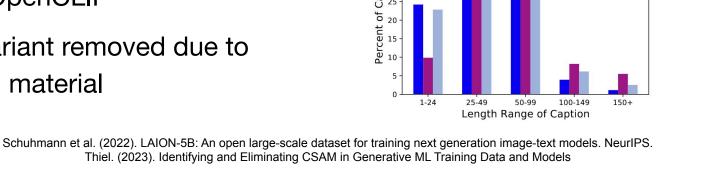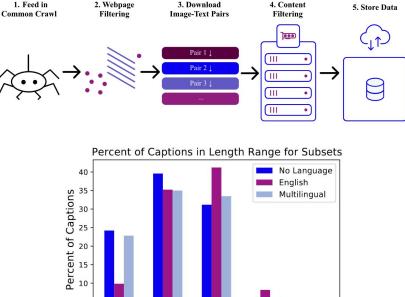
# LAION

- Used for pretraining

- Image and multilingual *raw* captions harvested from within Common Crawl

- Data behind Stable Diffusion and OpenCLIP

- 5B variant removed due to illegal material



Schuhmann et al. (2022). LAION-5B: An open large-scale dataset for training next generation image-text models. NeurIPS.
Thiel. (2023). Identifying and Eliminating CSAM in Generative ML Training Data and Models

# VQAv2

- Answer questions about images

- Task with multimodal inputs:
  - Image
  - Question

- Commonly tackled as classification but increasing efforts as NLG

- 1.1M image–question pairs with balanced distribution of answers



Goyal et al. (2017). Making the V in VQA matter: Elevating the role of image understanding in visual question answering. CVPR.

# NLVR2

- Binary classification task that requires jointly reasoning over a pair of images and a sentence.

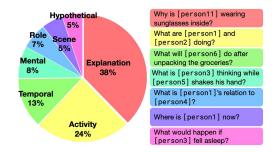- Human-created hard negatives.

- 107K examples in total.



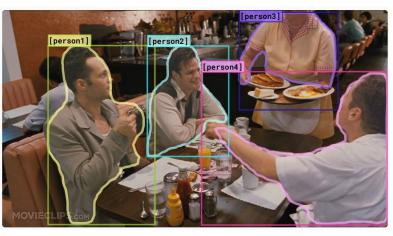*The left image contains twice the number of dogs as the right image, and at least two dogs in total are standing.*

Suhr et al. (2019). A Corpus for Reasoning about Natural Language Grounded in Photographs. ACL.

# Visual Commonsense Reasoning

- 290,000 multiple-choice VQA examples derived from movies.



- In addition to Question Answering, the dataset includes rationale selection too!



Zellers et al. (2019). From Recognition to Cognition: Visual Commonsense Reasoning. CVPR.

# Multi30K

- Multilingual aligned image–sentence dataset in many languages
  - English, German, French, Czech, Arabic, Japanese, Turkish, Ukranian

*A group of people are eating noodles.*

*Eine Gruppe von Leuten isst Nudeln.*

*Un groupe de gens mangent des nouilles.*

*Skupina lidí jedí nudle.*

Elliott et al. (2016). Multi30K: Multilingual English-German Image Descriptions. ACL Workshop.

# BOBSL

- BBC-Oxford British Sign Language Dataset

- Sign spotting and sentence localization tasks

- 1,400 hours of signed shows

  - Factual, entertainment, drama, comedy, children's shows

Albanie et al. (2021). BBC-Oxford British Sign Language Dataset. arXiv.

# Many Many More

- Visual Storytelling, e.g. VIST
- Grounded Referring Expression, e.g. Flickr30K Entities, Visual Genome
- Visual Entailment, e.g. SNLI-VE
- Vision & Language Navigation, e.g. RxR
- Visual Common Sense Reasoning: VCR
- Text-to-Image Generation, e.g. DALLEval
- Abstract reasoning, e.g. KiloGram, CRAFT
- Sign Language Processing, e.g. How2Sign

- *and more and more and more and more*

# Binding: Degree of Multimodality

- The content expressed in textual data depends on the purpose



Social media platforms often form 'echo chambers' that encourage users to only read content that confirms beliefs they already hold (Getty)

A woman in a grey suit is giving a speech

Weak ◄─────────────────────────► Strong

(Crawled)

(Crowdsourcing)

Rewriting crawled text improves performance on a variety of downstream multimodal tasks

Li et al. BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation. ICML 2022.
Nguyen et al. Improving multimodal datasets with image captioning. NeurIPS 2023.

# Ethical Issues

- Multimodal datasets are usually data scraped from the web with *unknown degrees of conformance*, or information about, licensing.
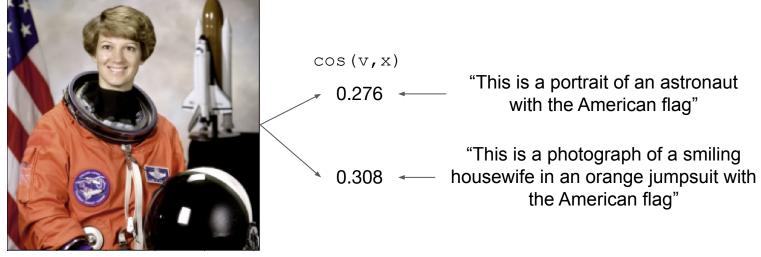


**CC BY:** This license allows reusers to distribute, remix, adapt, and build upon the material in any medium or format, so long as attribution is given to the creator. The license allows for commercial use.

- As of 2022, there are an estimated 2.5B CC-licensed objects online.

# The Problem with Scale

- Scale lets you build systems that perpetuate harmful stereotypes



(Eileen Collins, American astronaut)

$\cos(v,x)$

0.276 ← "This is a portrait of an astronaut with the American flag"

0.308 ← "This is a photograph of a smiling housewife in an orange jumpsuit with the American flag"

Birhane et al. (2021) Multimodal datasets: misogyny, pornography, and malignant stereotypes. arXiv.

# Q: How can we collect multimodal data that better reflects the diversity of the world?

# Motivation

**Languages**

- Mostly in English
- Or some Indo-European Languages



ENG: An unusual looking vehicle …
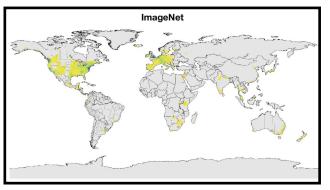
NLD: Een mobiel draaiorgel …

Example from van Miltenburg+ 2017

**Image sources**

- Mostly from ImageNet or COCO
- Reflecting North American and European cultures

**Implications for V&L models**

- Narrow linguistic/cultural domain
- No way to assess their real-world comprehension



Density map of geographical distribution of images in ImageNet (DeVries+, 2019)

F. Liu*, E. Bugliarello*, et al. Visually Grounded Reasoning across Languages and Cultures. EMNLP 2021.

# Typical Vision and Language (DATA)
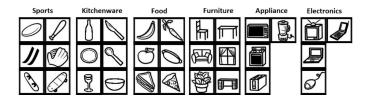
ImageNet (Deng et al. 2009)
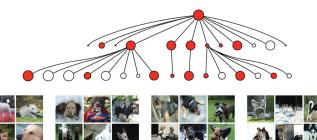
- Train visual encoders
- Millions of labelled images
- Derived from the WordNet hierarchy



canine → dog → working dog → husky

Common Objects in Context (Lin et al. 2014)

- Train and evaluate multimodal models
- 330K labelled images
  - 80 types of commonly occurring objects



Sports  Kitchenware  Food  Furniture  Appliance  Electronics

# Concrete Concepts in Cultural Context

- Some concepts are most immediately understood within a cultural background

  *Culture*: The way of life of a collective of people that distinguishes them from other people (Mora, 2013; Shweder et al. 2007).



**Pilota / Jai-alai**



**Sanxian / Shamisen**
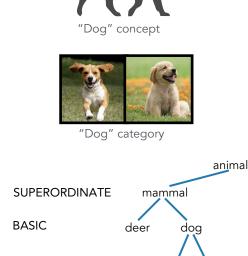


**Clavie**

# Concepts and Hierarchies

**Category:** objects with similar properties (Aristotle 40 BCE, …)

**Concept:** mental representation of a category (Rosch 1973)

Categories form a *hierarchy*

- Basic-level categories (Rosch 1976)

Somewhat universal

- Different cultures (Berlin 2014)
- Familiarity of individuals
  (Wisniewski and Murphy, 1989)



"Dog" concept



"Dog" category
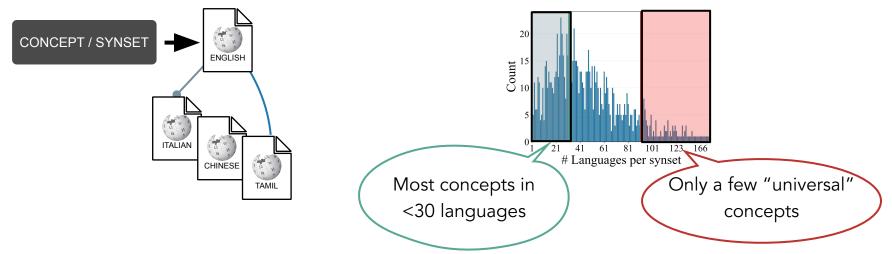
# Are ImageNet Concepts Cross-Lingual?

- ImageNet, COCO and Visual Genome use English WordNet concepts
- Question: estimate cross-linguality using Wikipedia as a proxy



CONCEPT / SYNSET → ENGLISH, ITALIAN, CHINESE, TAMIL

Most concepts in <30 languages

Only a few "universal" concepts

**MaRVL** Multicultural Reasoning over Vision and Language

📷 Representative of annotators' cultures

அ/A 5 typologically diverse languages
Independent, culture-specific annotations

MaRVL-id Bola basket    MaRVL-sw Mpira wa kikapu    MaRVL-tr Basketbol    MaRVL-zh 篮球    MaRVL-ta கூடைப்பந்தாட்டம்

# Visual Reasoning Task

- **Datapoint**: two images ($v_1$, $v_2$) paired with a sentence x

- **Task**: Predict whether x is a true description of the pair of images $v_1$ $v_2$



இரு படங்களில் ஒன்றில்
இரண்டிற்கும் மேற்பட்ட
மஞ்சள் சட்டை அணிந்த
வீரர்கள் காளையை அடக்கும்
பணியில் ஈடுப்பட்டிருப்பதை
காணமுடி.

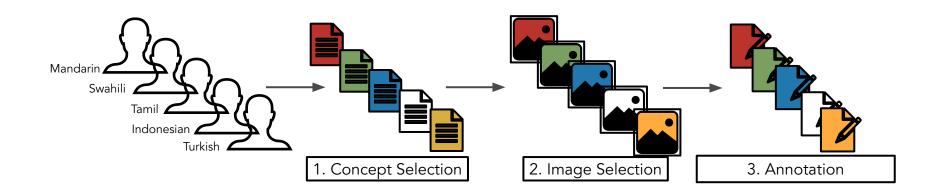True

x                                              y

Suhr et al. A Corpus for Reasoning about Natural Language Grounded in Photographs. ACL 2019.

# Collecting MaRVL data

Native speaker-driven protocol



Mandarin
Swahili
Tamil
Indonesian
Turkish

1. Concept Selection

2. Image Selection

3. Annotation

# **MaRVL** is created from Universal Concepts

- Taken from the *Intercontinental Dictionary Series* (Key & Comrie, 2015)
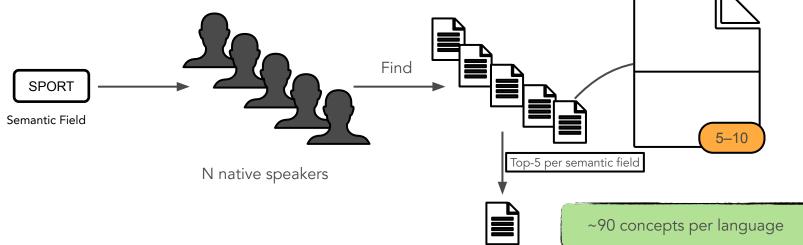  - 18/22 chapters with concrete objects & events

| Chapter | Semantic Field |
|---|---|
| Animal | Bird, mammal |
| Food and Beverages | Food, Beverages |
| Clothing and grooming | Clothing |
| The house | Interior, exterior |
| Agriculture and vegetation | Flower, fruit, vegetable, agriculture |
| Basic actions and technology | Utensil/tool |
| Motion | Sport |
| Time | Celebrations |
| Cognition | Education |
| Speech and language | Music (instruments), visual arts |
| Religion and belief | Religion |

⚠ OUR BIAS

38

# Step 1. Language-Specific Concepts

Defined by native speakers

- Commonly seen or representative in their culture
- Ideally, physical and concrete

SPORT

Semantic Field

Find

N native speakers

Top-5 per semantic field

5–10

~90 concepts per language

# Overview of Resulting Concepts

# Step 2. Image Collection

Collected by native speakers

- Representative of the language population
- NLVR2 (Suhr et al. ACL 2019) requirements

  1. Contains more than one instance of a concept

  2. Shows an instance of the concept interacting with other objects

  3. Shows an instance of the concept performing an activity

  4. Displays a set of diverse objects or features
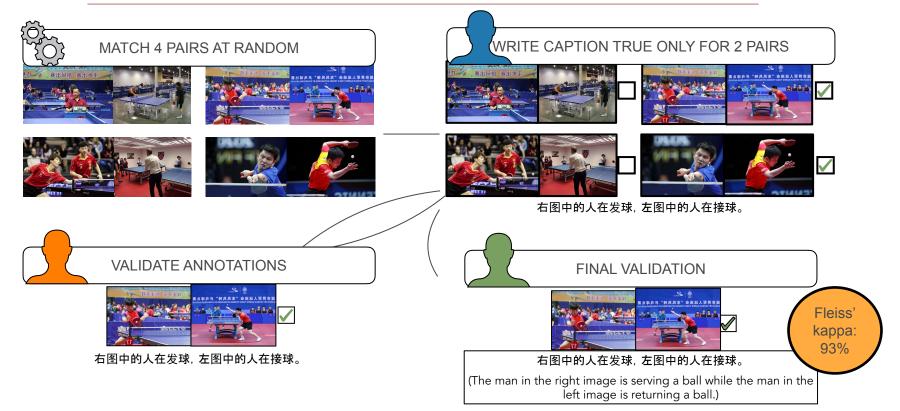


MaRVL-zh 花椰菜 (Cauliflower)



MaRVL-ta மோர் (Buttermilk)



MaRVL-sw Jembe (Shovel)



MaRVL-tr Rakı (Raki)

# Step 3. Language Annotation

MATCH 4 PAIRS AT RANDOM

WRITE CAPTION TRUE ONLY FOR 2 PAIRS

右图中的人在发球，左图中的人在接球。

VALIDATE ANNOTATIONS

右图中的人在发球，左图中的人在接球。

FINAL VALIDATION

右图中的人在发球，左图中的人在接球。

(The man in the right image is serving a ball while the man in the left image is returning a ball.)

Fleiss' kappa: 93%

42

# Dataset Examples



MaRVL-tr Kanun (çalgı)

Görsellerden birinde dizlerinde kanun bulunan birden çok insan var

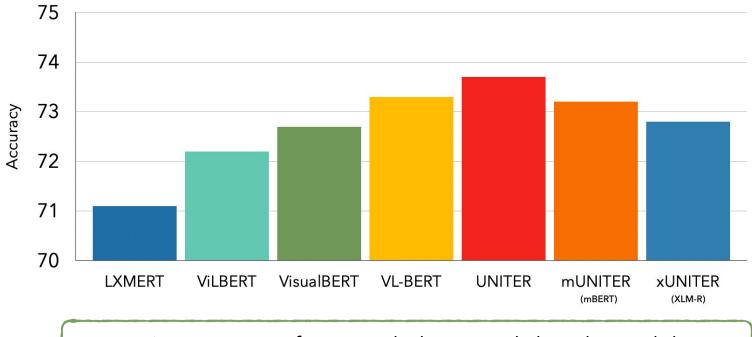(In one of the images, there are multiple people with qanuns on their knees)

Label: True



MaRVL-ta மை (Vada)

இரண்டு படங்களிலும் நிறைய மசால் வடைகள் உள்

(Both images contain a lot of masala vadas)
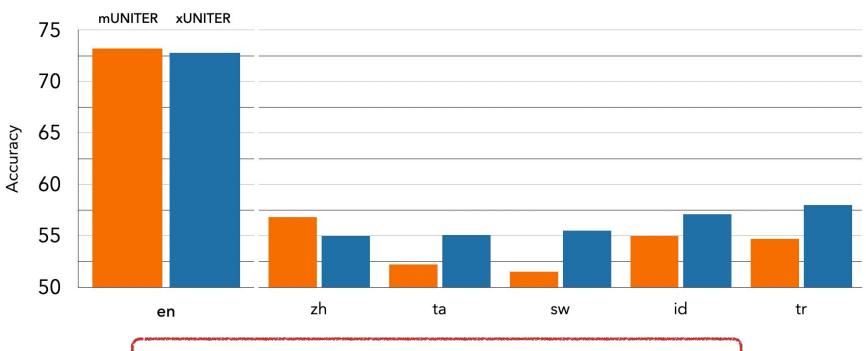
Label: False

43

# English NLVR2 Results (Sanity check)



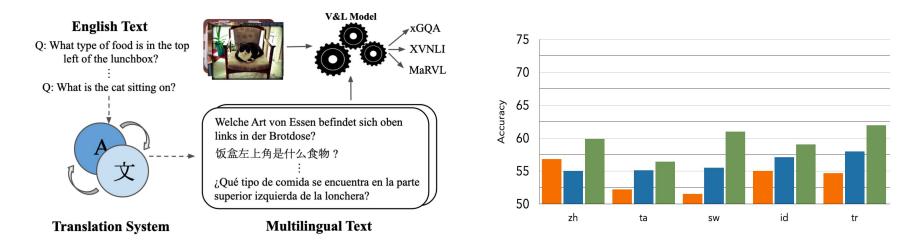m/xUNITER perform similarly to English-only models

# MaRVL Zero-shot Results



Zero-shot transfer: substantial drop in performance

45

# Pretraining with Translated Text

● Are the low zero-shot results caused by poor cross-lingual multimodal binding?



Cross-modal multilingual multimodal pretraining helps!
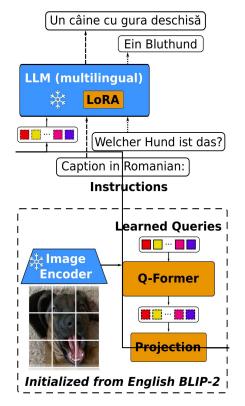
Qui et al. Multilingual Multimodal Learning with Machine Translated Text. Findings of EMNLP 2022.

# Supervised State of the Art: mBLIP

- Initialize from
  - BLIP-2
  - MT0-XL

- Use NLLB to pretrain on 96 languages
  - MSCOCO
  - CapFilt
  - VQAV2 & A-OKVQA
  - ImageNet as multilingua VQA



Geigle et al. mBLIP: Efficient Bootstrapping of Multilingual Vision-LLMs. 2023.

# Zero-shot State of the Art: PaliGemma

- Initialize from
  - Gemma 2B
  - SigLIP-So400m/14

- Pretrain on
  - Web Language Image (???)
  - CC3M-35L
  - VQ²A/VQG-CC3M-35L
  - OpenImages
  - Wikipedia Image Text

# Year-on-Year Improvements

- Clear benefit when using machine translated data

- Better visual encoders and language models can enable effective zero-shot transfer



A bar chart showing 5-Languages Accuracy:
- xUNITER (2021): 54.6
- TD-MML (2022): 61.6
- mBLIP (2023): 75.1
- PaliGemma (2024): 80.6

# 2. Data Representation

# Three Levels of Representation

- Perceptual
- Pre-processed features
- Raw input

❑ Yellow
❑ Has wheels
❑ Metal
❑ Five-door
❑ Can transport
❑ ...




Images: http://www.cs.columbia.edu/~vondrick/ihog/

# Perceptual Norms

- Ask people to write down the words that are triggered by textual stimuli.

- Stimuli: 541 noun concepts

- Norms are categorized into the likely knowledge source

| Moose | | |
|---|---|---|
| is large | 27 | visual–form and surface |
| has antlers | 23 | visual–form and surface |
| has legs | 14 | visual–form and surface |
| has four legs | 12 | visual–form and surface |
| has fur | 7 | visual–form and surface |
| has hair | 5 | visual–form and surface |
| has hooves | 5 | visual–form and surface |
| is brown | 10 | visual–color |
| hunted by people | 17 | function |
| eaten as meat | 5 | function |
| lives in woods | 14 | encyclopedic |
| lives in wilderness | 8 | encyclopedic |
| an animal | 17 | taxonomic |
| a mammal | 9 | taxonomic |
| an herbivore | 8 | taxonomic |

McRae et al. (2005). Semantic feature production norms for a large set of living and nonliving things. *Behavior research methods*, *37*(4), 547.

# Perceptual Norms: Pros / Cons

### Pros

- Seemingly simple task
- Rich features

### Cons

- Can it scale?
- Handling ambiguity

What does it mean to only understand symbols as defined by other symbols?



Searle. (1980). Minds, Brains and Programs. *Behavioral and Brain Sciences*, 3: 417–57

# Spatial and Pooled Visual Features

- Earliest work in neural-network era used pooled or spatial preserving features from a pretrained Convolutional Neural Network.



"where $CNN(I_b)$ transforms the pixels inside bounding box $I_b$ into 4096-dimensional activations of the fully connected layer immediately before the classifier." **Pooled features**

**Spatial features** "In our experiments we use the 14×14×512 feature map of the fourth convolutional layer before max pooling."

Karpathy & Fei-Fei (2015). Deep visual-semantic alignments for generating image descriptions. CVPR.
Xu et al. (2015). Show, attend and tell: Neural image caption generation with visual attention. ICML.

# Pre-processed Visual Features

- Faster R-CNN region-based feature vectors

  - Trained on the Visual Genome Dataset

  - The Region Proposal Network suggests the location of *regions of interest*.

  - RoI pooling performs spatial pooling in the final CNN layer to give a 2048D vector.



Our models are trained with 36 regions of interest extracted by a Faster R-CNN with a ResNet-101 backbone (Anderson et al., 2018).

Ren et al. (2015). Faster r-cnn: Towards real-time object detection with region proposal networks. *NeurIPS.*
Bugliarello et al. (2021). Multimodal Pretraining Unmasked: A Meta-Analysis and a Unified Framework of Vision-and-Language BERTs. TACL.

55

# Pre-processed: Pros / Cons

## Pros

- Long-established practice

- Usually an offline process:
do it once and forget

## Cons

- Large datasets require
specialized storage

- Not obvious how to
randomly augment data

- Specialist knowledge can be
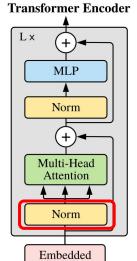opaque to newcomers

# Raw Input

- Directly process data from the raw images or speech signal.

- Images:
  - Vision Transformer (ViT)
  - Swin Transformer

- Speech
  - Spectrogram Transformer
  - AudioMAE

*Transformers | Davide Coccomini | 2021*

# Vision Transformer

- Good news! You are already almost an expert in how the Vision Transformer works

    - Split image into K patches

    - Embed each patch

    - Add position information

    - Encode using Transformer blocks that include an extra pre-norm layer for stability.



**Transformer Encoder**

L ×

MLP

Norm

Multi-Head Attention

Norm

Embedded Patches

Dosovitskiy et al. (2021). An image is worth 16x16 words: Transformers for image recognition at scale. ICLR.
Baevski & Auli (2018). Adaptive input representations for neural language modeling. *ICLR.*

58

# Extracting ViT Features

- Extract pooled features or patch-level features

To extract visual information from an image $x^i$, we use the visual encoder of a pre-trained CLIP [29] model. Next, we

For the CLIP encoders, we extract the feature grid before the pooling layers, resulting in an $N \times N$ grid, where $N = 7, 7, 12$ for the ViT-B/32, RN50x4 and RN50x16 variants of CLIP respectively.

Eichenberg et al. (2021). MAGMA--Multimodal Augmentation of Generative Models through Adapter-based Finetuning. *EMNLP*

# Raw input: Pros / Cons

## Pros

- Data augmentation is straightforward because you always have the raw input

- Fewer preprocessing steps means fewer creeping errors
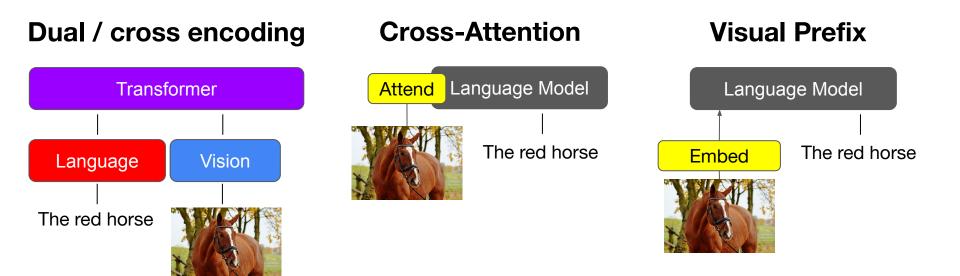
## Cons

- Smaller batches with an extra model on the GPU

- Potentially many inputs

# Summary

- Many options for how to represent your multimodal inputs

  - Language-oriented
  - Object / stuff oriented
  - Raw inputs

- **No universally best option** but raw inputs are promising because the visual encoding model can be fully differentiable
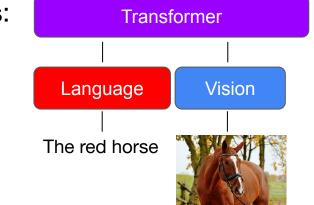
# 3. Modelling

# Main Approaches

**Dual / cross encoding**



**Cross-Attention**



**Visual Prefix**
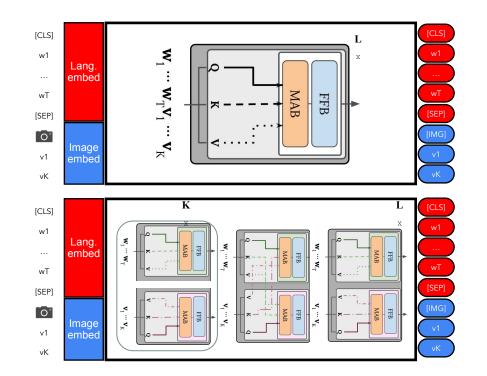
# Cross-encoding Models

- Emerged as a key modelling approach in 2019 with a flurry of approaches to creating visually-grounded BERT models.

- This is a form of *model-based fusion*

- The backbone consists of two components:
  - language encoder
  - visual encoder



Transformer

Language    Vision

The red horse

Tan & Bansal (2019). LXMERT: Learning Cross-Modality Encoder Representations from Transformers. EMNLP-IJCNLP.
Lu et al. (2019). ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. NeurIPS.
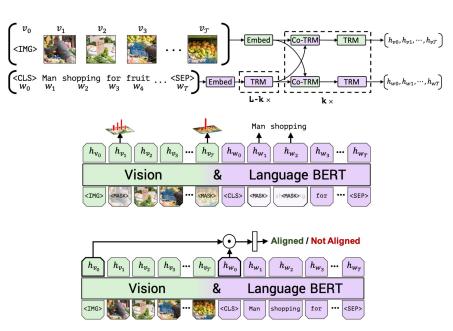
# Single- & Dual-Stream Architectures

- Single-stream
  - Concatenate inputs into one sequence

- Dual-stream
  - Process modalities independently
    - Intra-modal
    - Inter-modal



Bugliarello et al. 2021. Multimodal Pretraining Unmasked: A Meta-Analysis and a Unified Framework of Vision-and-Language BERTs. TACL.
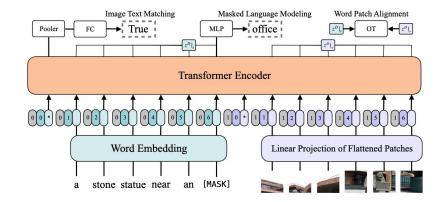
# ViLBERT

- Dual-stream model

- Initialized from BERT

- Visual features extracted from 10-36 regions using Faster-RCNN

- Pretrained on Conceptual Captions

  - Masked Language Modelling
  - Masked Region Classification
  - Image-Text Matching



Lu et al. (2019). ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. NeurIPS.
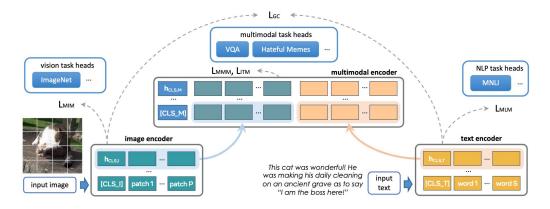
66

# ViLT

- Single-stream model

- Initialized from BERT

- Visual features extracted from ViT-B/32

- Pretrained on Conceptual Captions, Visual Genome, COCO, SBU Captions

    ○ Masked Language Modelling
    ○ Image-Text Matching
    ○ Word-Patch Alignment



Kim et al. (2021). ViLT: Vision-and-language transformer without convolution or region supervision. ICML.
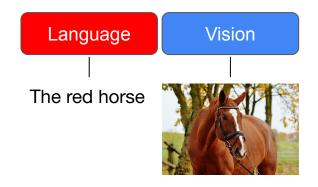
# FLAVA



- Dual-stream Visual features extracted from ViT-B/16

- Pretrained on PMD70M
    - Masked Language Modelling, Masking Image Modelling
    - Image-Text Matching, Masked Multimodal Modelling
    - Global Contrastive Matching

Singh et al. (2022). FLAVA: A foundational language and vision alignment model. CVPR.

# Dual-encoding Models

- Emerged as a sample-efficient alternative to cross-encoding.

- The backbone consists of two separate components:
  - language encoder
  - visual encoder



Radford et al. (2021). Learning transferable visual models from natural language supervision. ICML.
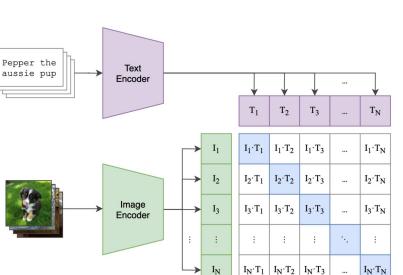
# CLIP

The red horse



- 12 Layer Transformer Encoder
- ViT or ResNet Visual Encoder
- Maximize the similarity of the embeddings of paired examples (I, T):

$$\mathcal{L}_{\text{InfoNCE}} = -\mathbb{E}\left[\log \frac{f(\mathbf{t}, \mathbf{i})}{\sum_{\mathbf{t}' \in T} f(\mathbf{t}', \mathbf{i})}\right]$$
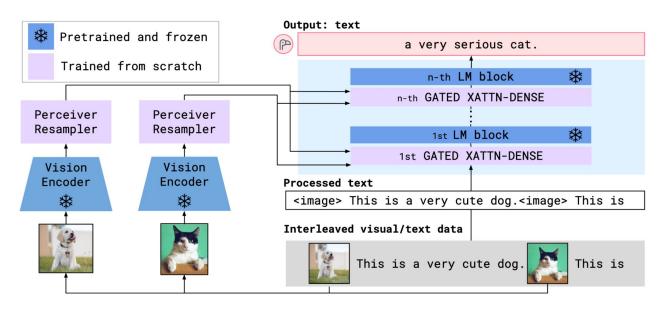
- Huge pretraining dataset of unclear provenance



Radford et al. (2021). Learning transferable visual models from natural language supervision. ICML.

70

# Cross-Attention

Attend Language Model

The red horse



Alayrac et al. NeurIPS 2022. Flamingo: a Visual Language Model for Few-Shot Learning

71

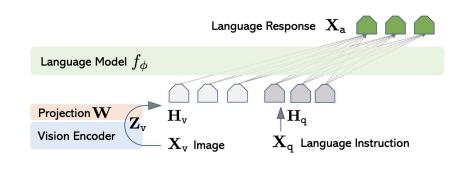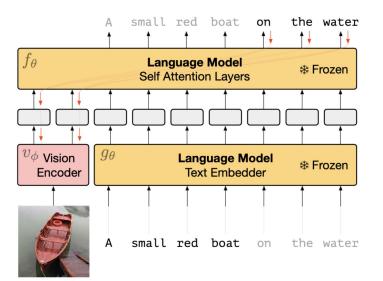# Visual Prefix

- Exploit the representations learned during large-scale modality specific pretraining



Tsimpoukelli et al. NeurIPS 2021. Multimodal Few-Shot Learning with Frozen Language Models.
Liu et al. NeurIPS 2023. Visual Instruction Tuning.

# Current Vision and Language Models
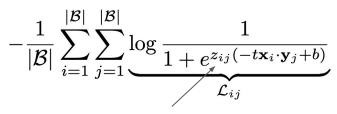
- Current models mostly follow this blueprint:

    1. Choose pretrained modality-specific components
    2. Learnable bridge between those components
    3. Dataset to estimate the parameters in the bridge
    4. Multi-stage finetuning strategy

# Modality-Specific Components

|  | **Vision Encoder** | **Language Model** |
|---|---|---|
| **LLAVA** | CLIP ViT-L/14 | Vicuna-13B |
| **Qwen-VL** | OpenCLIP ViT-bigG | Qwen-7B |
| **MM1** | ViT-L | 1.3B LLM |
| **PaliGemma** | SigLIP-So400M/14 | Gemma-2B |

https://ai.google.dev/gemma/docs/paligemma
https://llava-vl.github.io/
https://github.com/QwenLM/Qwen-VL

74

# SigLIP Image Encoder

"Unlike standard contrastive learning with softmax normalization, the sigmoid loss operates solely on image-text pairs and does not require a global view of the pairwise similarities for normalization."

$$-\frac{1}{|\mathcal{B}|}\sum_{i=1}^{|\mathcal{B}|}\sum_{j=1}^{|\mathcal{B}|}\log\underbrace{\frac{1}{1+e^{z_{ij}(-t\mathbf{x}_i\cdot\mathbf{y}_j+b)}}}_{\mathcal{L}_{ij}}$$

Label of the image-text pair: 1 if matched else -1



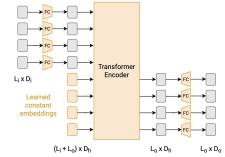Zhai et al Sigmoid Loss for Language Image Pre-Training. ICCV 2023.
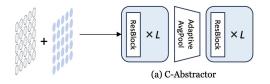
# Learnable Bridge

- LLAVA and PaliGemma: Use a single linear layer to map the image output embeddings into the language model word embedding space.

- Qwen-VL:

**Position-aware Vision-Language Adapter**: To alleviate the efficiency issues arising from long image feature sequences, Qwen-VL introduces a vision-language adapter that compresses the image features. This adapter comprises a single-layer cross-attention module initialized randomly. The module uses a group of trainable vectors (Embeddings) as query vectors and the image features from the visual encoder as keys for cross-attention operations. This mechanism compresses the visual feature sequence to a fixed length of 256. The ablation about the number of queries is shown in Appendix E.2. Additionally, considering the significance



- MM1: Convolutional-Abstractor
  - ResNet Block followed by an Adaptive Pooler



(a) C-Abstractor

Manas et al. MAPL: Parameter-efficient adaptation of unimodal pre-trained models for vision-language few-shot prompting. EACL 2024.
Cha et al. Honeybee: Locality-enhanced Projector for Multimodal LLM. CVPR 2024
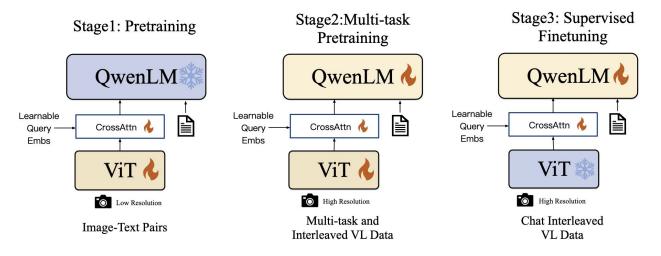
# Training Dataset

- LLAVA: 595K image–caption examples filtered from CC3M

  Qwen-VL          1.4 billion examples (77% English / 23% Chinese)

  MM1              2+ billion mixture of image–text examples

  PaliGemma        1 billion mixture of multilingual image caption, VQA, and in-the-wild datasets

- The larger models are pretrained on **in-house data**
  - PaliGemma: WebLI (1B+), Qwen-VL (220M), MM1 (1B+)

# Data Processing

- Encode the text using the language model tokenizer
- Encode the image using the image encoding model

- Image-position embeddings for multi-image sequences

- PaliGemma-specific
  - Location co-ordinate tokens
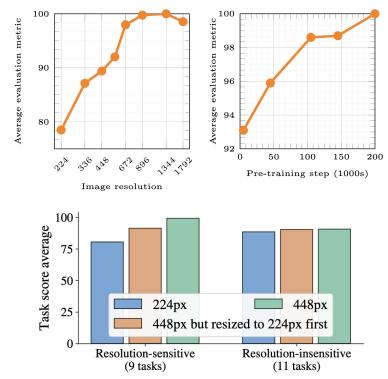  - Segmentation tokens

# Training Strategy

- Qwen-VL, PaliGemma, and MM1 use multi-stage training strategies with different types of data and different image resolutions

# Open Questions

- How quickly will we realize these benefits in smaller models?

- Do LMMs really need 1 billion examples to learn a bridge?

- What happens to performance when we develop new tasks that involve weaker visual–linguistic bindings?

# Summary

- Cross-encoding:

    - Many advances in which parts of the input contribute to loss

    - Shift from regions-of-interest to Vision Transformers

- Dual-encoding:

    - Excellent cross-domain transfer to a wide range of problems

- Visual Prefix Learning:

    - Exploit the benefits of single-modality pretraining

# Q: Does an image captioning model need to learn everything in-weights?
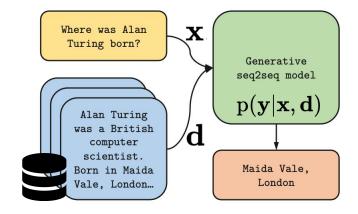
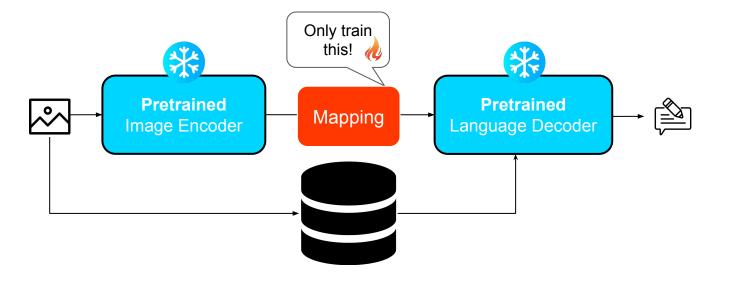# Retrieval Augmented Generation

- Combine the power of in-weights learning with in-context adaptation through retrieval augmentation

- Given a datastore of facts, knowledge, documents, etc.

  - Combine the most relevant items from the datastore (d) with the input (x) for your task



Izacard and Grave. EACL 2021. Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering.

# Motivation

- Main trend in V&L is training bigger models on more data
- Alternative is emerging that re-uses independent backbone models
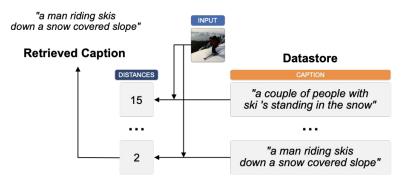- Can we further improve performance with retrieval augmentation?

# PAELLA Model

# Retrieval System

- Build a FAISS datastore: store high-dimensional vectors
  - Captions of images represented with CLIP embeddings

- Retrieve k nearest-neighbours captions from datastore
  - Image embedding compared against datastore caption vectors



Johnson et al. (2019). Billion-scale similarity search with GPUs. IEEE Big Data

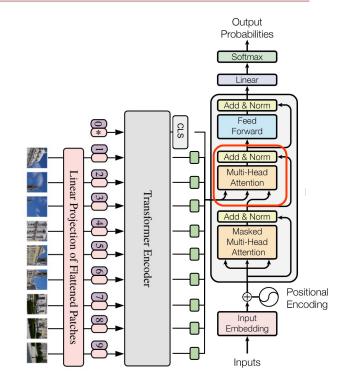# Trained Cross-Attention Layers

- We insert a randomly initialized **cross-attention mechanism** to attend to the visual encoder output embeddings

| Rank | Params |
|------|--------|
| d=128 | 553M |
| **d=8** | 34M |

$$\mathrm{Att}(\mathbf{X}\mathbf{W}^Q, \mathbf{X}\mathbf{W}^K, \mathbf{X}\mathbf{W}^V)$$

$$\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V \in \mathbb{R}^{\mathrm{d\_enc} \times d}$$

# Experimental Protocol

- Encoder: Multilingual CLIP
- Decoder: XGLM-2.9B

- Training data:
    - 566K captions sampled uniformly from COCO-35

- Evaluation: XM-3600
    - 3600 geographically-diverse images
    - 36 languages: 100 captions per image
    - 5 low-resource languages (L5):
        - Bengali, Cusco Quechua,
          Maori, Swahili, Telugu



Examples images from XM3600

# Results

| | Data | Trained $\theta$ | L36 | L5 |
|---|---|---|---|---|
| PaLI | 12B | 17B | 53.6 | - |
| Lg $_{COCO-35}$ | 19M | 2.6B | 15.0 | 12.5 |
| mBLIP: BLOOMZ-7B | 135M | 800M | 23.4 | 6.7 |
| BB+CC$_{COCO-35 + CC-35}$ | 135M | 800M | 28.5 | 22.4 |
| mBLIP: mT0-XL | 489M | 124M | 28.3 | 7.9 |
| **PAELLA** | **566K** | **30M** | 26.2 | 20.7 |

PAELLA is competitive against models with 35-863x more training data, and 4-87x more trained parameters
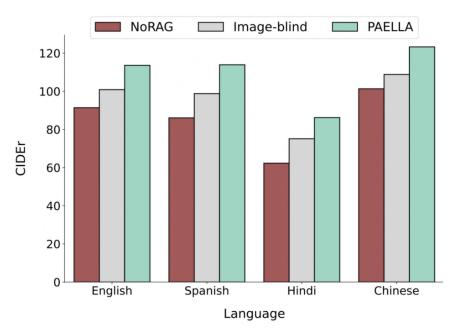
# Zero-shot Multilingual Transfer

- **PAELLA**<sub>MONO</sub> is a variant trained on 566K examples in English COCO

- Outperforms **Lg** trained on 19.8M examples in the machine translated COCO-35 dataset

| | Data | Trained θ | L36 | L5 |
|---|---|---|---|---|
| Lg: Thapliyal et al. $_{COCO-35}$ | 19M | 2.6B | 15.0 | 12.5 |
| **PAELLA**$_{MONO}$ | 566K$_{en}$ | 30M | 15.5 | 12.1 |

# Value of Retrieval Augmentation

Consistent improvements from multilingual retrieval augmentation across the core languages in the XM3600 evaluation data

# Qualitative Example



类似图片显示:

ऐसी ही तस्वीरें दिखाती हैं:

Imágenes similares muestran:

Similar images show:

*the owl is perched outside in front of the people*
*an owl sitting a top a table during the daytime*
*an owl is sitting on a perch at a camp site*
*the fuzzy owl is sitting on a tree branch*

A caption I can generate to describe this image in
**english** is:

**PAELLA**

en: *"an owl sitting on top of a tree"*

es: *"un búho sentado en una rama de un árbol"*
(an owl sitting on a tree branch)

hi: *"एक उल्लू एक पेड़ की टहनी पर बैठा है"*
(an owl is sitting on a tree branch)

zh: *"一只 猫头鹰 站在 树上"*
(an owl standing in a tree)

**NoRAG**

en: *"a large black and white picture of a bird"*

es: *"un pájaro posado en la parte superior de un edificio"*
(a bird perched on the top of a building)

hi: *"एक पेड़ के पास खड़ा एक पक्षी"*
(a bird standing near a tree)

zh: *"一只 长颈鹿 坐在 树枝 上"*
(a giraffe sitting on a branch)

93

# Try it yourself (in English)



Demo: https://huggingface.co/spaces/RitaParadaRamos/SmallCapDemo

# Why Compositionality?

- Given recent advances in MLLMs, we should work on tasks that require more sophisticated logical or commonsense reasoning

    - RecipeQA (Yagcioglu et al. 2018)

    - ScienceQA (Lu et al. 2022)

- **Sequential Multimodal Compositional Generalization** requires models to reason across a sequence of related multimodal inputs

RecipeQA: A Challenge Dataset for Multimodal Comprehension of Cooking Recipes. Yagcioglu et al. EMNLP 2018.
Learn to explain: Multimodal reasoning via thought chains for science question answering. Lu et al. NeurIPS 2022.
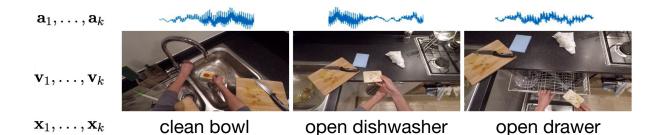
# EPIC Kitchens 100

- Head-mounted recordings from 45 different kitchens

- Rich annotations of objects

- Simple action descriptions

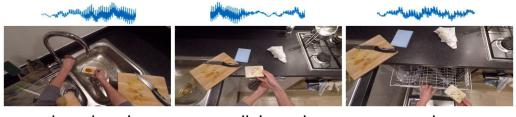  - 93 verb classes

  - 300 object classes



Rescaling Egocentric Vision: Collection, Pipeline and Challenges for EPIC-KITCHENS-100. Damen et al. IJCV 2022.

# CompAct Dataset

- Multimodal sequences consisting of aligned segments
  - Audio recording (not speech)
  - Video frames
  - Short description

$\mathbf{a}_1, \ldots, \mathbf{a}_k$

$\mathbf{v}_1, \ldots, \mathbf{v}_k$

$\mathbf{x}_1, \ldots, \mathbf{x}_k$

clean bowl       open dishwasher       open drawer

# Compositional Generalization Tasks
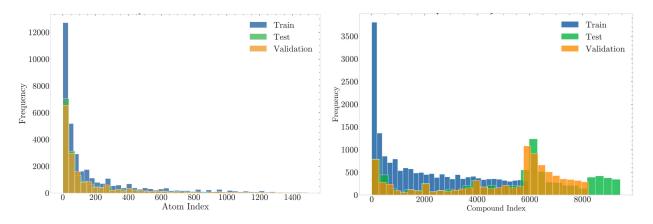


clean bowl      open dishwasher      open drawer

1. Next Description Prediction using a language model

$$p_\theta(\mathbf{y} = \mathbf{x}_{k+1} | \mathbf{x}_{1:k}, \mathbf{v}_{1:k}, \mathbf{a}_{1:k})$$

2. Atom Classification
   - Verb     $p_\theta(\mathbf{y} = v | \mathbf{x}_{1:k}, \mathbf{v}_{1:k}, \mathbf{a}_{1:k})$
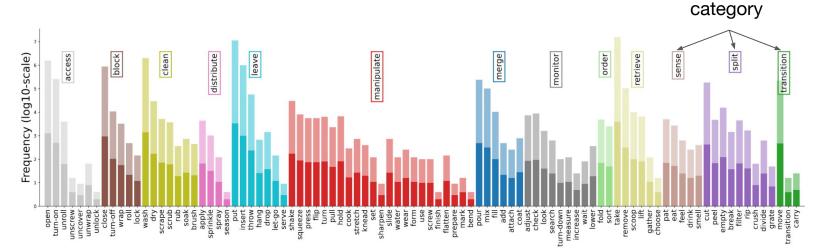   - Object   $p_\theta(\mathbf{y} = o | \mathbf{x}_{1:k}, \mathbf{v}_{1:k}, \mathbf{a}_{1:k})$

# Atoms and Compositionality

- We use Maximum Compound Divergence (Keysers et al. 2020) to create a dataset that requires compositional generalization
  - Noun and verbs extracted from simple descriptions



Measuring compositional generalization: A comprehensive method on realistic data. Keysers et al. ICLR 2020.

# CompAct Dataset Statistics

- 8,766 multimodal sequences
  - 50% training, 25% validation, 25% test



Distribution of verb classes in CompAct

# Models

## Baselines

Trained on CompAct:

- **T**ext-only
- **V**ision & **L**anguage
- **O**bject & **L**anguage
- **A**udio & **L**anguage
- **V**ision, **A**udio, & **L**anguage
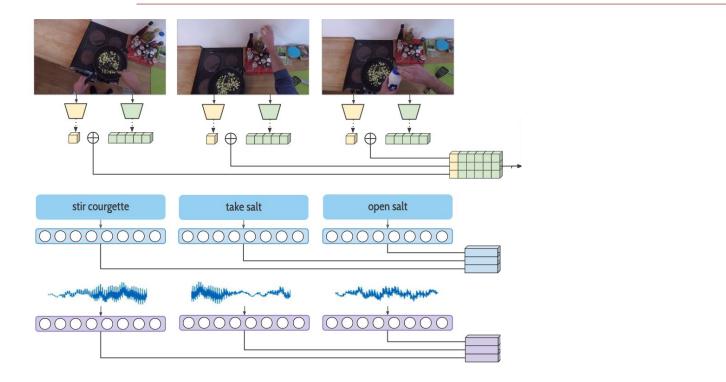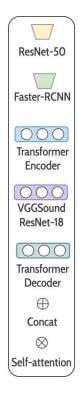- **O**bject, **A**udio, & **L**anguage

## Pretrained (M)LLMs

*k*=8-shot prompting:

- LLaMa-2 Chat 7B
- IDEFICS-9B
- OpenFlamingo-9B
- Otter-7B

No guarantee that these models need to compositionally generalize

# Baseline Architecture

# Example Prompt Templates

## LLaMA-2 Prompt `(k=3)`

```
Predict the next narration given 3 sequential previous narrations from a cooking video
put down bowl . move frying pan . pick up spatula => put down spatula
move yoghurt . put down bowl . pick up yogurt => put yoghurt
put down bowl . grab wok . move tap => lather wok
pick up tins . put down tins . move bowl =>
```

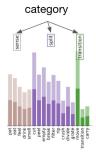## IDEFICS Prompt `(k=1)`

```
Predict the next action narration given 3 sequential previous actions (image-narration pairs) in a cooking video.
put down bowl <Image 1> . move frying pan <Image 2> . pick up spatula <Image 3> => put down spatula
pick up tins <Image 1> . put down tins <Image 2> . move bowl <Image 3> =>
```
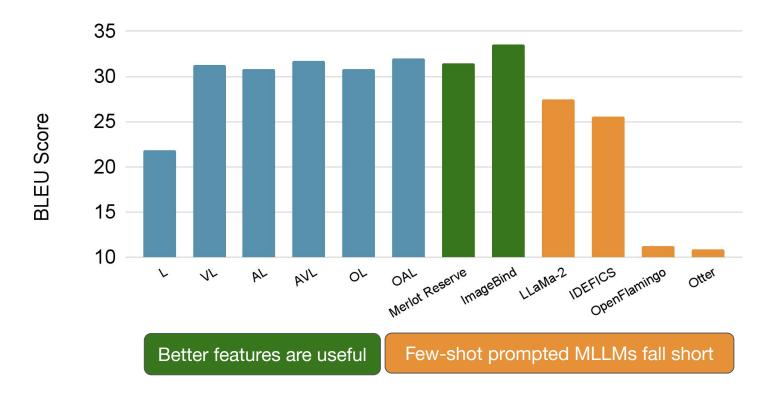
# Experimental Details

- Video segments represented by one keyframe
  - Based on detected objects from Faster R-CNN

- Baselines tokenize multi-part tokens as one token
  - olive oil -> olive_oil

- MLLMs use own tokenizers

- All experiments used one 16GB NVIDIA V100.
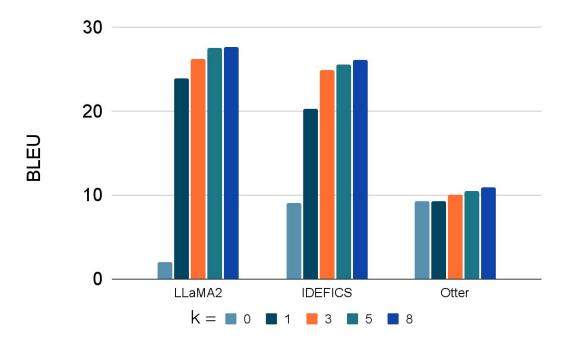  Each baseline can be trained in < 1 hour

# Evaluation Measures

- BLEU score

- Exact Match
  - How often does the model predict exactly the expected noun or verb?

- Categorical Accuracy
  - Does the model predict a noun or verb in the same semantic category?
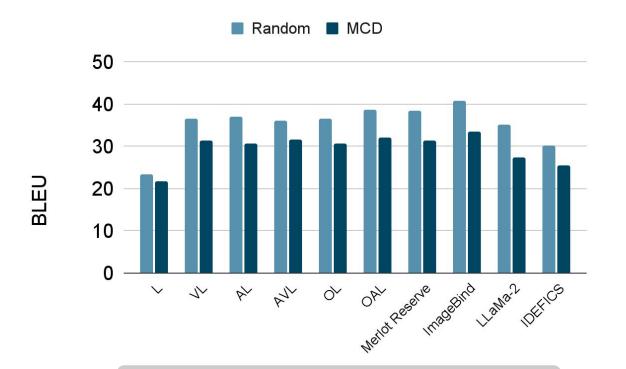
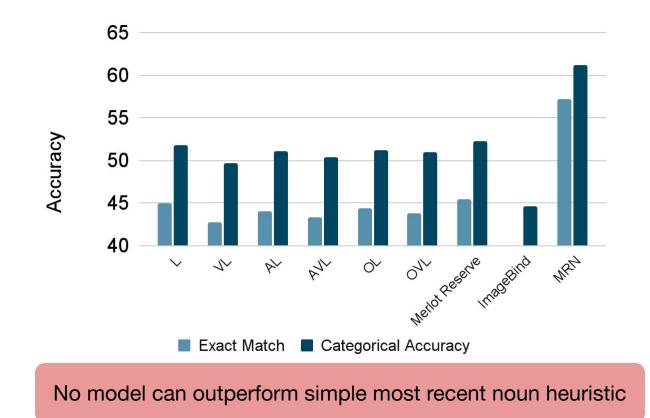# Next Utterance Prediction

# Few-shot Prompting



Clear benefit from increasing the number of in-context examples

# Random I.I.D. split data is easier



All models generalize better to randomly split data

# Noun Prediction



No model can outperform simple most recent noun heuristic

# Qualitative Examples


clean bowl


open dishwasher


open drawer

L: close dishwasher
AL: close dishwasher
VL: close dishwasher
AVL: dry bowl
IB: put bowl in dishwasher
LLaMA-2: get clean
IDEFICS: put bowl away


put pan in drainer


pick up mug


pick up sponge

L: put sponge
AL: sponge mug
VL: sponge mug
AVL: sponge mug
IB: sponge mug
LLaMA-2: put sponge in sink
IDEFICS: sponge mug

# Conclusions

- Sequential Multimodal Compositional Generalization is challenging new task where better multimodal features improve performance

- We find no evidence that ICL or large-scale multimodal pretraining can solve this task

- Future work includes
  - integrating even better features into the baseline
  - fine-tuning MLLMs using LoRA
  - including more keyframes in the visual input

# 4. Understanding Multimodal Models
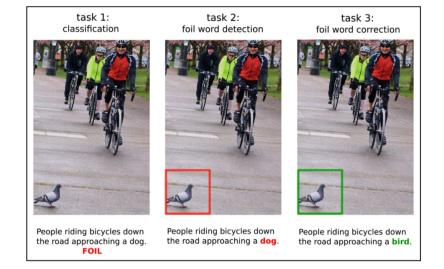
# Beyond Benchmarking

- Many questions about what drives the success of these models?
    - Better contextualization: make better use of the multimodal inputs
    - Acquire certain "skills", e.g. counting or localization
    - Understand linguistic structures
    - Something else?

- Model-internal behaviour
    - Attention mechanism patterns

- Probing
    - Tasks related to different skills

# FOIL Captions

- Do V&L models really understand the relationship between words and images?

- Crowdsource datasets that contain contextually plausible but incorrect image–text pairs



task 1: classification

task 2: foil word detection

task 3: foil word correction

People riding bicycles down the road approaching a dog. **FOIL**

People riding bicycles down the road approaching a **dog**.

People riding bicycles down the road approaching a **bird**.

Shekhar et al. (2017). FOIL it! Find One mismatch between Image and Language caption. ACL.
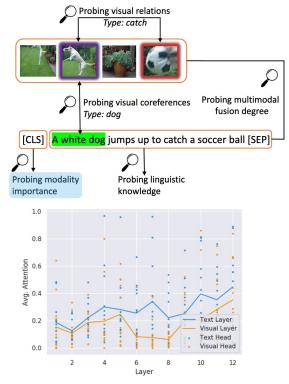
# Vision and Language Understanding Evaluation

- Suite of five model probing tasks

- **Modality Influence**: Estimate the layer-wise contribution of each modality to the `[CLS]` embedding:

$$I_{M,j} = \sum_{i \in S} \mathbb{1}(i \in M) \cdot \alpha_{ij}$$

  - The UNITER model relies more on textual features when fusing modalities throughout the model



Cao et al. (2020). Behind the scene: Revealing the secrets of pre-trained vision-and-language models. ECCV.

# VALSE Benchmark

- Test visio-linguistic capabilities with image-sentence foil pairs

- Image-sentence matching task

  - Existential quantifiers
  - Semantic number
  - Counting
  - Prepositional relations
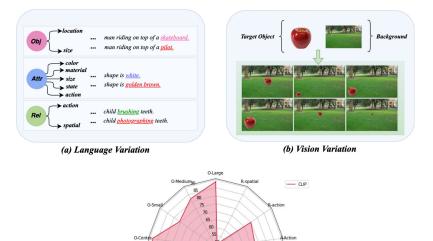  - Action replacement / swap
  - Co-reference

A small copper vase with some flowers / exactly one flower in it.

| Metric | Model | Avg. |
|--------|-------|------|
| | Random | 50.0 |
| | GPT1* | 60.7 |
| | GPT2* | 60.1 |
| | CLIP | 64.0 |
| $acc_r$ | LXMERT | 59.6 |
| | ViLBERT | 63.7 |
| | 12-in-1 | **75.1** |
| | VisualBERT | 46.4 |

$$p(caption, img) > p(foil, img)$$

Parcalabescu et al. (2022). VALSE: A Task-Independent Benchmark for Vision and Language Models Centered on Linguistic Phenomena. ACL.

# VL-CheckList

- Evaluate V&L models based on automatic manipulations to vision and language data.

- Image-Sentence matching task

- Radar chart overviews based on object / attribute / relationship variations


(a) Language Variation


(b) Vision Variation



Zhao et al. (2022). An Explainable Toolbox for Evaluating Pre-trained Vision-Language Models. EMNLP.

# Subject-Verb-Object Probes

- Large-scale dataset with SVO triplets mined from Conceptual Captions and 14K images and with crowdsourced captions

- Foil detection formulation



*Children* cross the street.

Pos — child, cross, street
Neg — lady, cross, street

A *animal* lays in the grass.

Pos — animal, lay, grass
Neg — woman, lay, grass

Hendricks and Nematzadeh. (2021). Probing Image–Language Transformers for Verb Understanding. ACL.

# WinoGround

- 1,600 text-image pairs to evaluate compositional understanding
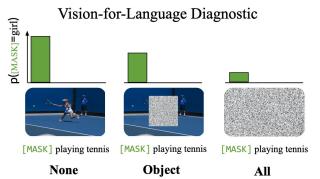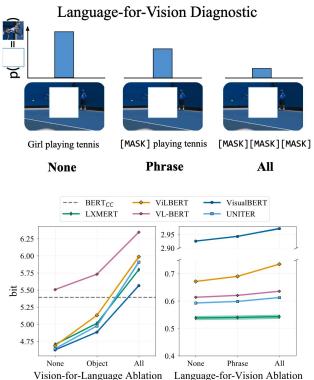


some plants surrounding a lightbulb



a lightbulb surrounding some plants

- Images sourced **with permission** from Getty.

- Differences are categorised into: swap dependent, swap-independent, and visual differences

Thrush et al. (2022). Winoground: Probing vision and language models for visio-linguistic compositionality. CVPR.

# Vision-for-Language?



Vision-for-Language Diagnostic

Language-for-Vision Diagnostic

p([MASK]=girl)

[MASK] playing tennis — **None**
[MASK] playing tennis — **Object**
[MASK] playing tennis — **All**

Girl playing tennis — **None**
[MASK] playing tennis — **Phrase**
[MASK][MASK][MASK] — **All**

- Pair of diagnostic evaluations that can be applied to any model that makes MLM and MRC predictions.



BERT$_{CC}$   ViLBERT   VisualBERT
LXMERT   VL-BERT   UNITER

bit

Vision-for-Language Ablation   Language-for-Vision Ablation

None   Object   All      None   Phrase   All

Frank et al. (2021). Vision-and-Language or Vision-for-Language? On Cross-Modal Influence in Multimodal Transformers. EMNLP.

# Summary

- Understanding and analysis is a vibrant area of research

- Foil detection is the most popular methodology

- Witnessing a methodological shift
  - attention analyses → linguistically-informed analyses
  - hand-crafted datasets
  - simpler accuracy-based metrics

# 5. Future Directions

# Physical Understanding

- Predicting and explaining physical actions in the world will become of increasing importance as we create embodied agents
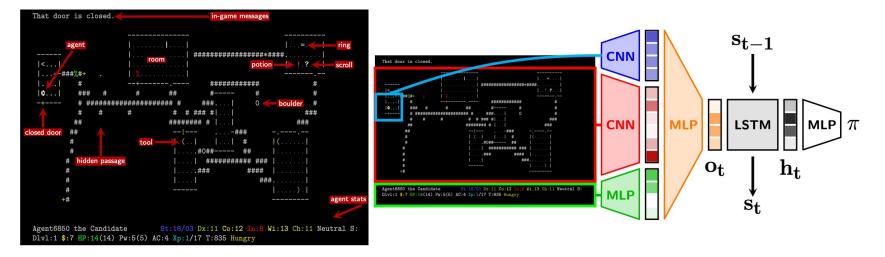


Q: How many objects are prevented by the tiny green triangle from falling into the basket?

Q: What is the color of the last object that collided with the tiny red circle?

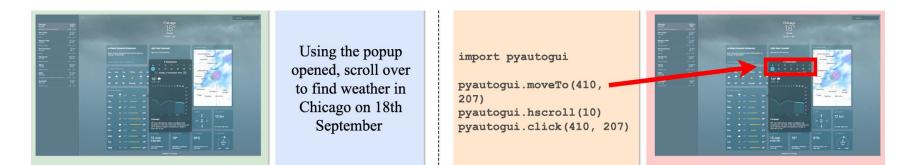Q: If any of the other objects are removed, will the tiny green circle end up in the basket?

Ates et al. (2022). CRAFT: A Benchmark for Causal Reasoning About Forces and inTeractions. ACL.

# Text-based Video Games

- Learning to act in procedurally-generated video game environments with rich contexts, action spaces, and long-term rewards



Küttler et al. (2020). The NetHack Learning Environment. NeurIPS.
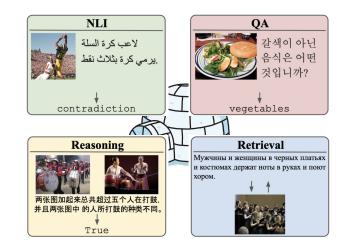
# Multimodal Interaction

- OmniAct combines multimodal understand with program execution to solve a variety of tasks that humans perform with their computers
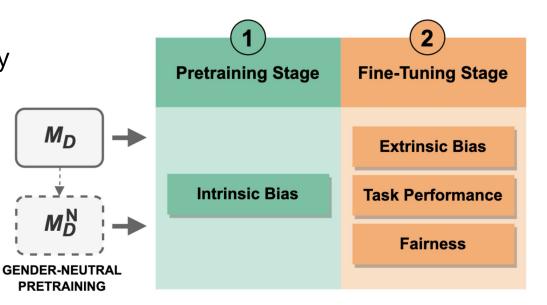


Kapoor et al. ECCV 2024. A Dataset and Benchmark for Enabling Multimodal Generalist Autonomous Agents for Desktop and Web.

# Multilinguality

- The majority of Vision and Language research is in English

- We need resources, models, and evaluations to create useful multilingual multimodal models

- High-quality data requires:
  - time
  - money
  - community engagement



Bugliarello et al. ICML 2022. IGLUE: Image-Grounded Language Understanding Evaluation.

# Bias and Fairness

- What are the intrinsic biases learned during multimodal pretraining and how do they affect downstream task performance?
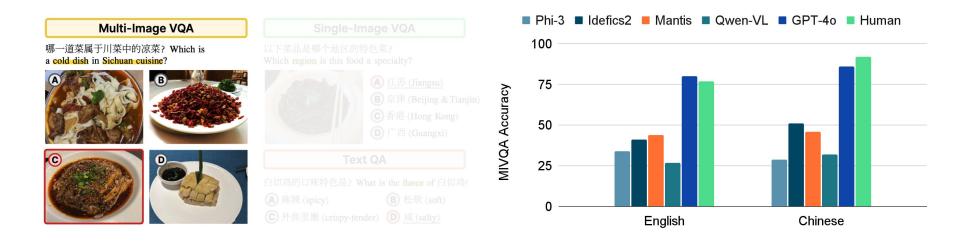


Cabello et al. EMNLP 2023. Evaluating Bias and Fairness in Gender-Neutral Pretrained Vision-and-Language Models

*Next*

# Fine-grained Multimodal Data …

- **Domain-specific fine-grained VQA data**

  - Chinese food
  - Highly-detailed taxonomy
  - Human questions
  - Three version of the task
  - Private data



Li et al. FoodieQA: A Multimodal Dataset for Fine-Grained Understanding of Chinese Food Culture. 2024.

# is challenging for LMMs



**Multi-Image VQA**

哪一道菜属于川菜中的凉菜？Which is a **cold dish** in Sichuan cuisine?

(A) (B)
(C) (D)

Single-Image VQA

以下菜品是哪个地区的特色菜？
Which **region** is this food a specialty?

(A) 江苏 (Jiangsu)
(B) 京津 (Beijing & Tianjin)
(C) 香港 (Hong Kong)
(D) 广西 (Guangxi)

Text QA

白切鸡的口味特色是？What is the **flavor** of 白切鸡？

(A) 麻辣 (spicy)    (B) 松软 (soft)
(C) 外焦里嫩 (crispy-tender)    (D) 咸 (salty)

Legend: ■ Phi-3  ■ Idefics2  ■ Mantis  ■ Qwen-VL  ■ GPT-4o  ■ Human

MIVQA Accuracy chart (English, Chinese)

API-based model barely outperforms a naive English annotator

Huge gap to fill between API-based and open-weights models

130

# *wonders:* Community-driven Multimodal Data

- Large-scale data collection with crowdworkers: most budgets cannot scale

**Core team**
- Define data & task
- Recruit and coordinate
  Researchers (eg PhDs)

**Researchers**
- Coordinate annotation
- Verify annotations

**Crowdworkers**
- Create annotations

- Our bet: **People care about their culture**



WIKIPEDIA

VISIPEDIA
Images, segments
annotations, links,
GUIs, diagnostics

Users

Experts

Image databases
flickr
IMAGENET

Annotators
amazon mechanical turk
Artificial Artificial Intelligence

Automata

Vision
scientists

Created with mapchart.net

# Gamified Data Collection

**Category** 􀆅

**Concept** 􀆅



Manden skærer en flæskesteg til aftensmad

Submit

💻 | En mand med en gris | 👍 / 👎

# Q: What if we treated language as vision?

# NLP in the Era of Scale

Treviso et al. 2023. Efficient Methods for Natural Language Processing: A Survey. TACL
Wei et al. 2022. Emergent Abilities of Large Language Models. TMLR

# NLP for **All** Written Languages

- There are 3,000 written languages
  - 400 with >1M speakers

- NLP usually covers 100 languages
  - Technological exclusion for billions

- We need systems for all languages, not just those that are high-resource

(Data and groupings from Joshi et al. ACL 2020)



| Class | Speakers |
|-------|----------|
| 0 | 1.2B |
| 1 | 30M |
| 2 | 5.7M |
| 3 | 1.8B |
| 4 | 2.2B |
| 5 | 2.5B |

van Esch et al. Writing System and Speaker Metadata for 2,800+ Language Varieties. LREC 2022.
Joshi et al. The State and Fate of Linguistic Diversity and Inclusion in the NLP World. ACL 2020.

How can we create high-quality Natural Language Processing tools for all written languages?

# NLP is a pipeline ...

Raw text — **Kierkegaard (d. 1855) was a Golden Age philosopher**

Separate punctuation

Normalize — **Kierkegaard ( d . 1855 ) was a Golden Age philosopher**

Split into sub-words

Sub-word Encoding
Character / UTF-32

Tokenize — **Kier #keg #aard ( d . 1885 ⋯ philosopher**

Each token is a
real-valued vector

Embedding
Lookup

| -0.1 | 0.2 | -0.4 | -0.1 | 0.3 | -0.2 | -0.1 | | 0.5 |
| 0.5 | 0.9 | 0.8 | 0.5 | 0.5 | 0.3 | 0.5 | ... | 0.7 |
| ... | ... | ... | ... | ... | ... | ... | | ... |
| 0.3 | -0.1 | 0.8 | 0.3 | 0.9 | 0.7 | 0.3 | | 0.2 |

Neural Network
(GPT, LLaMA, …)

Model — Model

**Syntactic / Semantic analysis**

# … that is easily broken

**Søren** Kierkegaard (d. 1855) was a Golden Age philosopher

ኢትዮጵያ አፍሪካ ውስጥ ናት

Not found

Tokenize

**UNK** **Kier** **#keg** **#aard** **(**
**d** **.** **1885** ⋯ **philosopher**

UNK  UNK  UNK  UNK

Embedding Lookup

| 0.0 | 0.2 | -0.4 | -0.1 | | 0.5 |
| 0.0 | 0.9 | 0.8 | 0.5 | ⋯ | 0.7 |
| … | … | … | … | | … |
| 0.0 | -0.1 | 0.8 | 0.3 | | 0.2 |

| 0.0 | 0.0 | 0.0 | 0.0 |
| 0.0 | 0.0 | 0.0 | 0.0 |
| … | … | … | … |
| 0.0 | 0.0 | 0.0 | 0.0 |

Model

Model

Model

This issue disproportionately affects low-resource languages

# The Vocabulary Bottleneck

- NLP is an **open vocabulary problem** and the

  ability of a model is determined by its vocabulary:

  1. tokens, characters, sub-words, etc.

- This creates a bottleneck in two places:

  1. *Representational bottleneck* in the Embedding layer
  2. *Computational bottleneck* in the Output layer

# Where's the sweet spot?

Long sequence lengths

*Bytes* (ByT5)

(CANINE, CharFormer) *Characters*

*Small subword vocab* (BERT: 30K)

(RoBERTa: 50K) *Medium subword vocab*

*Full unicode coverage* (ByT5, CANINE)

(XML-R: 250K) *Large subword vocab*

*Words* (word2vec: 3M)

Large vocabulary

# Main idea: treat language as vision



Raw text

Normalize

Tokenize

Embedding Lookup

Model

Render text as an image

ኢትዮጵያ አፍሪካ ውስጥናት

Randomly mask

Model

# The PIXEL Model



$$\frac{1}{M} \sum_{m \in M} (Y_m - \hat{Y}_m | \mathbf{x}_{\searrow m})^2$$

Transformer Decoder — 8 Layers

Transformer Encoder — 12 Layers

3 CLS Embedding & Span Mask $m$ patches

2 Projection + Position Embedding

16pixel x 16pixel patch

Google Noto Fonts

PyGame / PangoCairo

1 Render Text as Image

My cat ⊂ᴎᴎᴑ enjoys eating warm oatmeal for lunch and dinner.

# Rendered Text is Compact

- PIXEL encoding produces sequence lengths that are at least as long as as BERT.

  ○ Universal Dependencies datasets with human reference segmentations

  ○ No length penalty for any language, unlike some LLMs (Ahia et al. 2023)



Proportion of text that is encoded as $k$ subwords / patches.

# Pretraining

- **English Dataset**: English Wikipedia and Books Corpus
- **Masking:** 25% Span Masking
- **Maximum sequence length**: 529 patches (`16x8464` pixels)
- **Compute**: 8 x 40GB A100 GPUs for 8 days
- **Parameters**: 86M encoder + 26M decoder

There is only 0.05% non-English text in our pretraining data (estimated by Blevins and Zettlemoyer 2022)

The **Great Wall of China** (traditional Chinese: 萬里長城; simplified Chinese: 万里长城; pinyin: *Wànlǐ Chángchéng*)

Pretrained model: https://huggingface.co/Team-PIXEL/pixel-base

# A new type of generative model



Try for yourself

100K steps         500K steps         1M steps

# Downstream Tasks

- **Datasets**: Universal Dependencies, MasakhaNER, GLUE, Zeroé

- **Models:**

| | Parameters | Pretraining Data |
|---|---|---|
| PIXEL$_{BASE}$ | 86M | English Wikipedia + Bookcorpus |
| BERT$_{BASE}$ | 110M | — |
| CANINE-C | 127M | 104-languages from Wikipedia |

Similar pretraining setup

Tries to solve the same problem using UTF-32

Code: https://github.com/xplip/pixel

# Dependency Parsing Results



Legend: ■ BERT  ■ PIXEL

| | ENG | ARA | COP | HIN | JPN | KOR | TAM | VIE | ZHO |
|---|---|---|---|---|---|---|---|---|---|
| BERT UNK | 0% | 1% | 94% | 33% | 46% | 85% | 82% | 5% | 73% |

PIXEL vastly outperforms BERT on unseen scripts

# Named Entity Recognition



Emir of **Kano** turban **Zhang** wey don spend **18 years** for **Nigeria**

BERT — CANINE — PIXEL

F1

PIXEL outperforms BERT on the non-Latin script

PIXEL outperforms the multilingually pretrained CANINE-C

# Text Rendering Matters

- The original text renderer produces many nearly-identical patches
    - This is representation- and compute-wasteful



(a) Continuous rendering (CONTINUOUS):

(b) Structured rendering (BIGRAMS):

(c) Structured rendering (MONO):

(d) Structured rendering (WORDS):

# Sentence-level Tasks: GLUE



Legend: ■ BERT ■ PIXEL ■ PIXEL-BIGRAM

Chart showing GLUE task scores (0-100) for MNLI, QQP, QNLI, SST-2, COLA, STS-B, MRPC, RTE, WNLI.

Bigram text rendering produces better models

# Scaling Down ↓

- Better text rendering can create effective models at smaller scales

# Application: Historical Document Processing

- Worldwide efforts to digitize historic documents (Groesen 2015)

- Typical pipeline for enabling access is:

  a. Scan documents into high-quality digital formats
  b. Perform OCR on those documents (one-off process)
  c. Search through documents using OCR annotations

### What if we could do this without OCR?

# Caribbean Newspapers, 1718–1876

- Collaboration with researchers that are interested in tracking newspapers notices about escaped slaves

  - What was the given name?

  - What reward was offered?

  - Who was the contact person?

- Dataset of 1.65M scanned pages

# PIXEL for Historical Documents

- Historical document-aware Pretraining

  - Mixture of scanned newspapers and synthetic newspaper-like text generated from Wikipedia and Bookcorpus datasets

  - All input data is scaled to 368x368 and split into 16x16 patches

# Visual Question Answering in Newspapers

- Frame this as a Visual Question Answering Task

  - Render the question
  - Render the clipping on a canvas
  - Annotate context of answer

- Train the model to predict the label of the answer

How much reward is offered?

WHEREAS a Molatto Boy (his Name Dench) belonging to a young Man lately arrived from the East-Indies, absented himself on Monday the 20th Instant. He had on when he went away a Thicket Frock and Waistcoat, Leather Breeches, and a blue Surtuit Coat, with a red Collar.

Any Person that will apprehend the abovementioned Boy, or give any Intelligence where he may be taken, shall receive a Reward of Three Guineas. He is about five Feet high, with short black Hair, not of the woolly Kind.

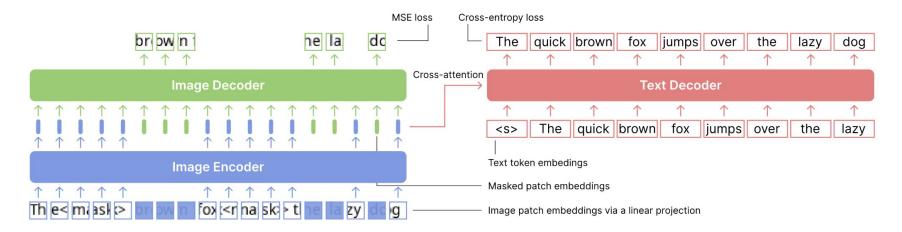N. B. If taken, to be brought to the Sign of the George, Queen-Ann-Street, Cavendish-Square.
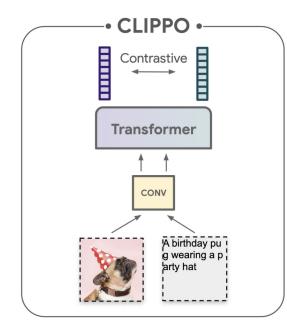
# Results



Surprisingly good performance compared to a model trained on manually transcribed text

# Patch and Text Prediction

- Combine patch and token prediction

# Joint Multimodal Reasoning



Tschannen et al. CVPR 2023. CLIPPO: Image-and-Language Understanding from Pixels Only.

# Open Questions

- Interpretability:
  - Does this work based on orthographic similarity or is it learning grammatical representations of text from pixels?

- Multilinguality and scale:
  - How should we train a multilingual PIXEL encoder?
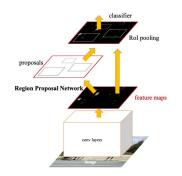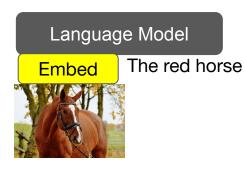    - Language-based or script-based data selection

# Wrap-up

# 1. Datasets

some sheep walking in the middle of a road
a herd of sheep with green markings walking down the road
a herd of sheep walking down a street next to a lush green grass covered hillside.
sheared sheep on roadway taken from vehicle, with green hillside in background.
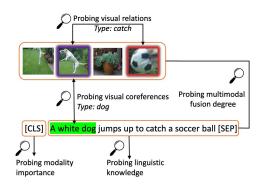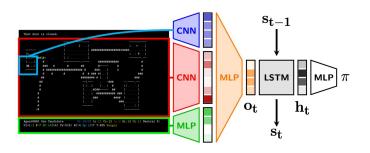a flock of freshly sheared sheep in the road.

# 2. Representation

classifier

RoI pooling

proposals

Region Proposal Network

feature maps

conv layers

image

# 3. Modelling

Language Model

Embed    The red horse

Probing visual relations
Type: catch

Probing visual coreferences
Type: dog

Probing multimodal
fusion degree

[CLS] A white dog jumps up to catch a soccer ball [SEP]

Probing modality
importance

Probing linguistic
knowledge

# 4. Understanding

That door is closed.

CNN

CNN

MLP

MLP

LSTM

MLP  π

$s_{t-1}$

$o_t$

$h_t$

$s_t$

Agent6850 the Candidate   St:18/03 Dx:11 Co:12 In:8 Wi:13 Ch:11 Neutral S:
Dlvl:1 $:7 HP:16(16) Pw:5(5) AC:4 Xp:1/17 T:635 Hungry

# 5. New Directions